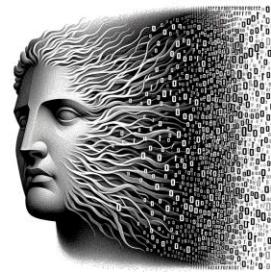




ASPLOS 2024



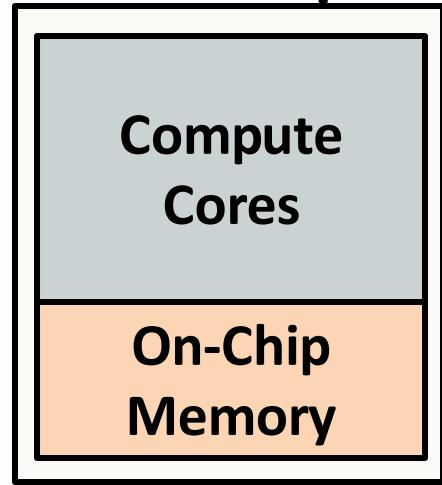
Atalanta:

A Bit is Worth a “Thousand” Tensor Values

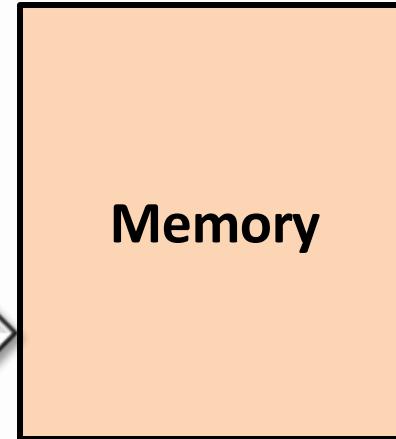
Alberto Delmas Lascorz, Mostafa Mahmoud, Ali Hadi Zadeh, Miloš Nikolić,
Kareem Ibrahim, **Christina Giannoula**, Ameer Abdelhadi,
and Andreas Moshovos^{1,2}



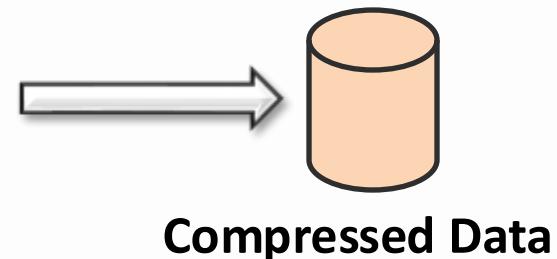
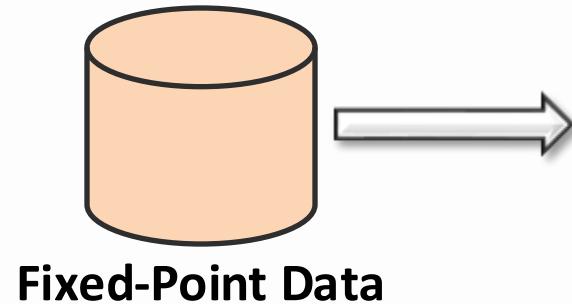
On-Chip



Off-Chip



Energy: ~100x
Latency: ~50x

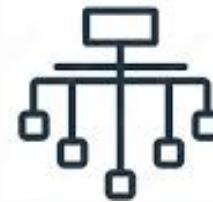




Inference & Training



Transformers



Classification



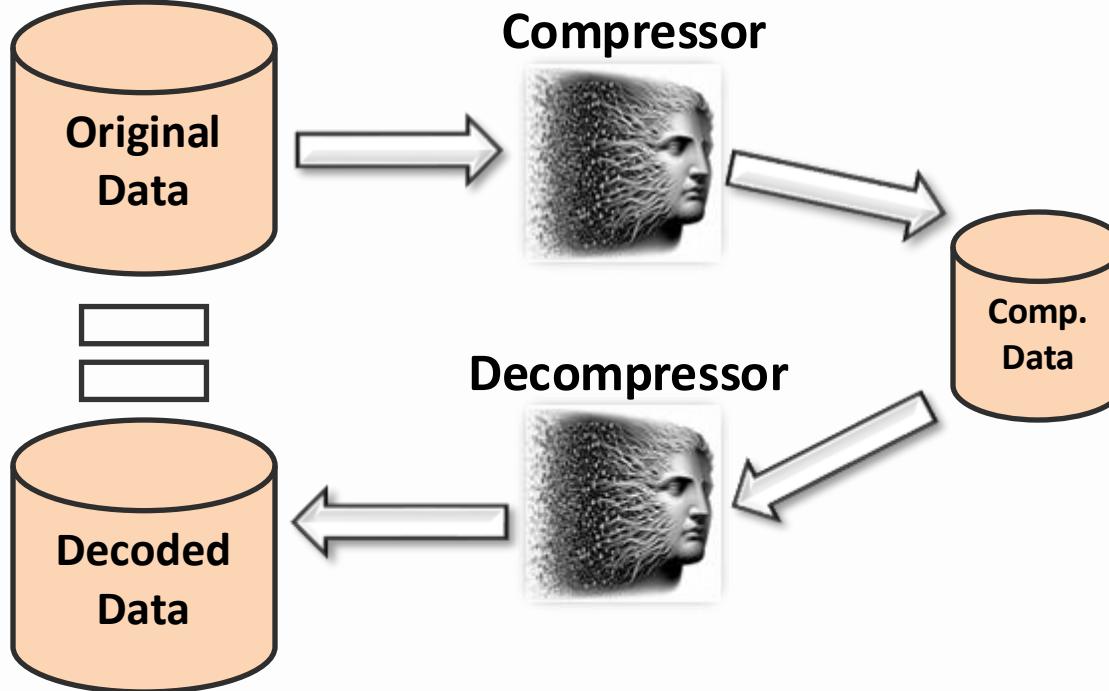
Recommendation



Detection

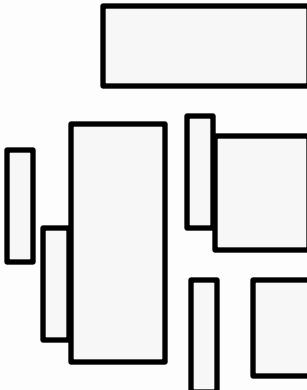


Lossless

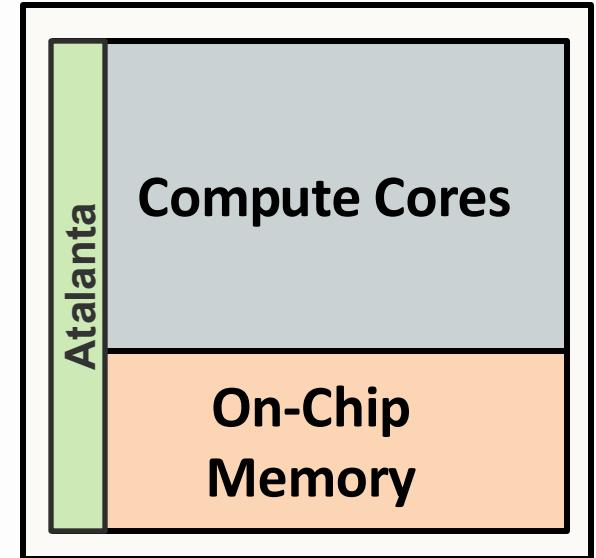




Lossless



Transparent

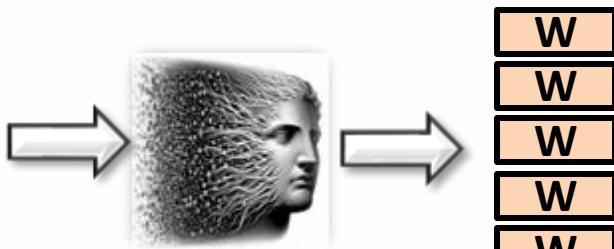




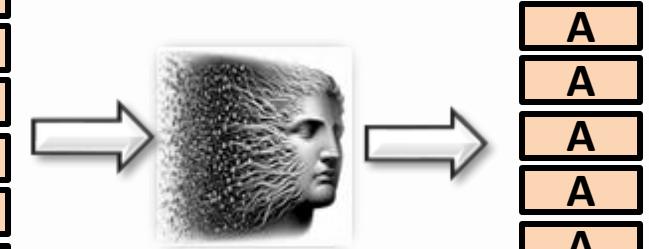
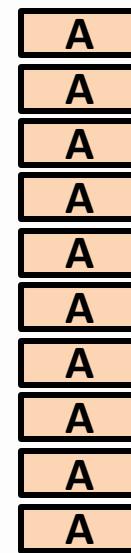
Lossless

Transparent

High Compression



60%



48%

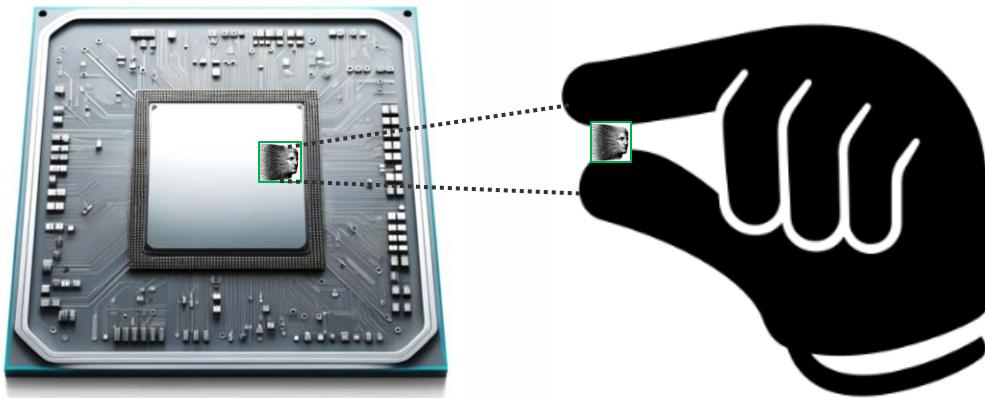


Lossless

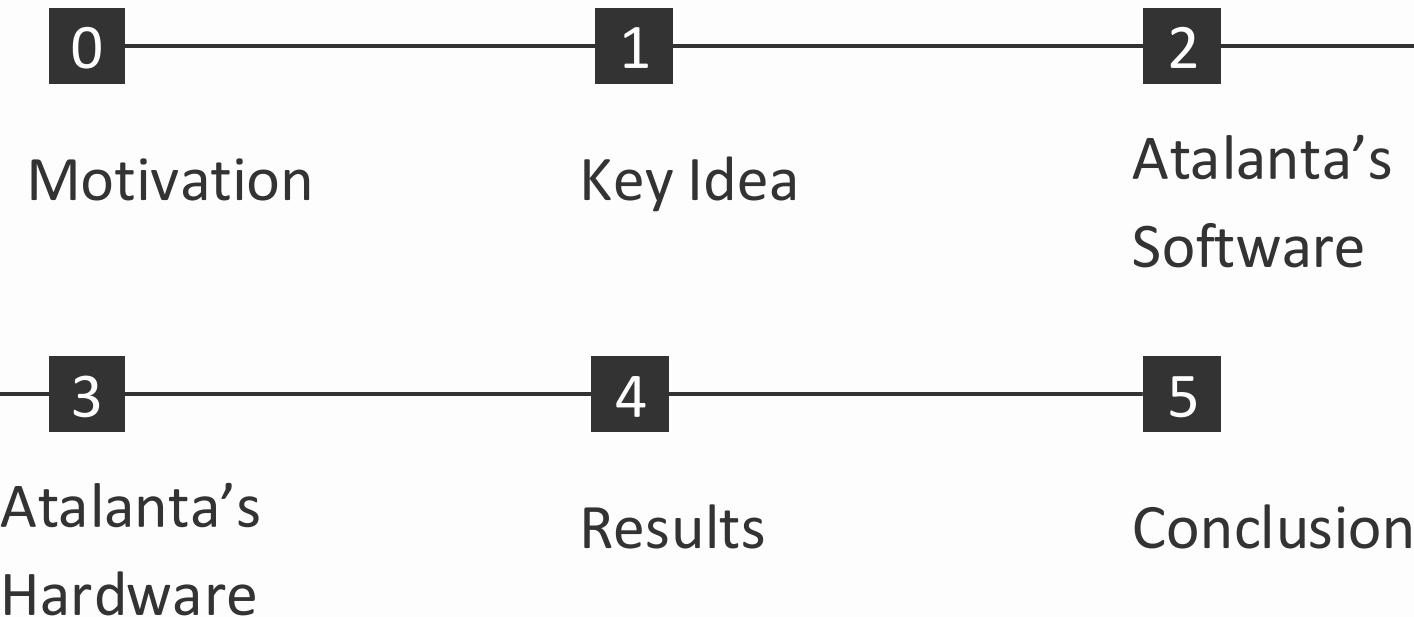
Transparent

High Compression

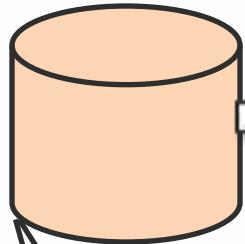
Low Cost



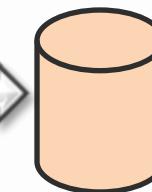
Roadmap



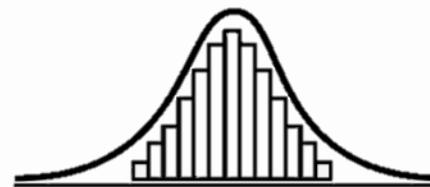
Fixed-Point Data



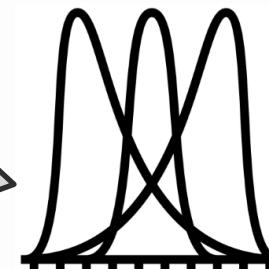
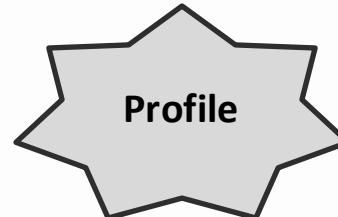
Compressed Data

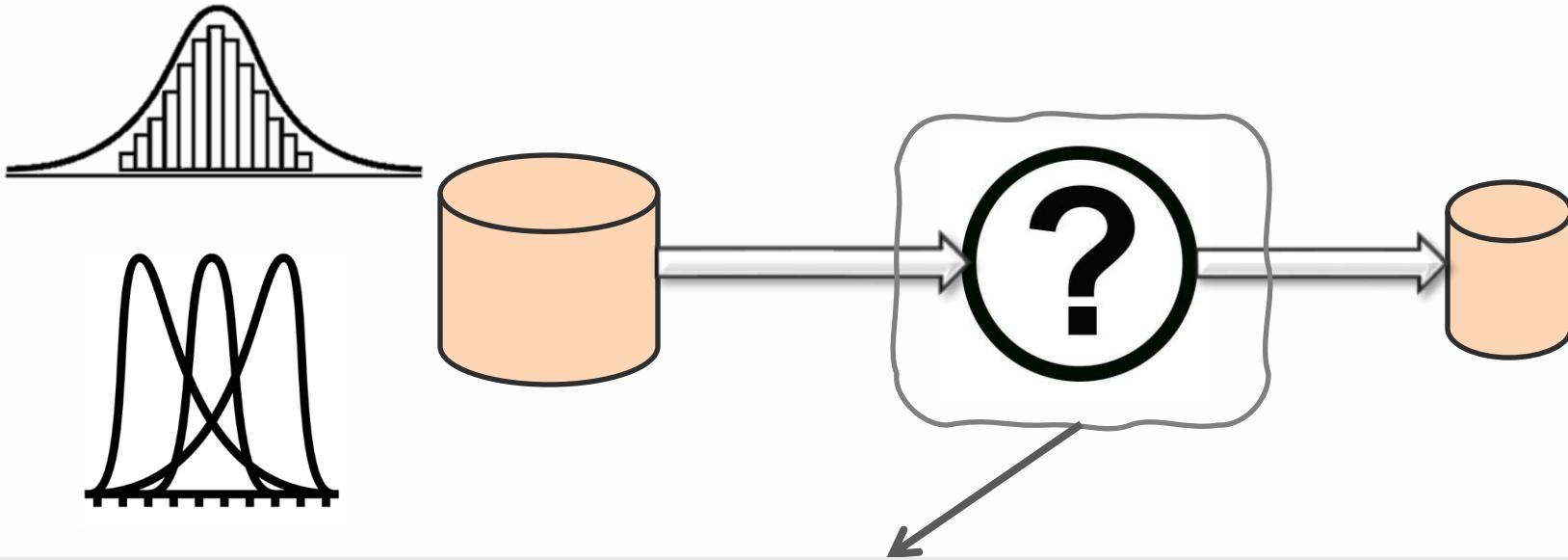


Weights



Activations





Arithmetic Coding?

Arithmetic Coding

INT8
Tensor

12, 0, 23, 45, 67, 127, 18, 22, 88, 103, 234, 22, 1, 0, 2, 3, 5, 8, 19, 9, 0, 9, 8, 20, 28, 220, 20, 20,
1, 0, 19, 9, 0, 9, 8, 20, 28, 220, 20, 20, 244, 223, 2, 1, 1, 0, 1, 0, 12, 0, 23, 45, 67, 127, 18, 22,
.....
2, 3, 5, 8, 28, 220, 20, 20, 244, 223, 2, 1, 1, 0, 1, 0, 234, 22, 1, 0, 2, 3, 19, 9, 0, 67, 127, 18, 22,
88, 103, 5, 8, 9, 8, 20

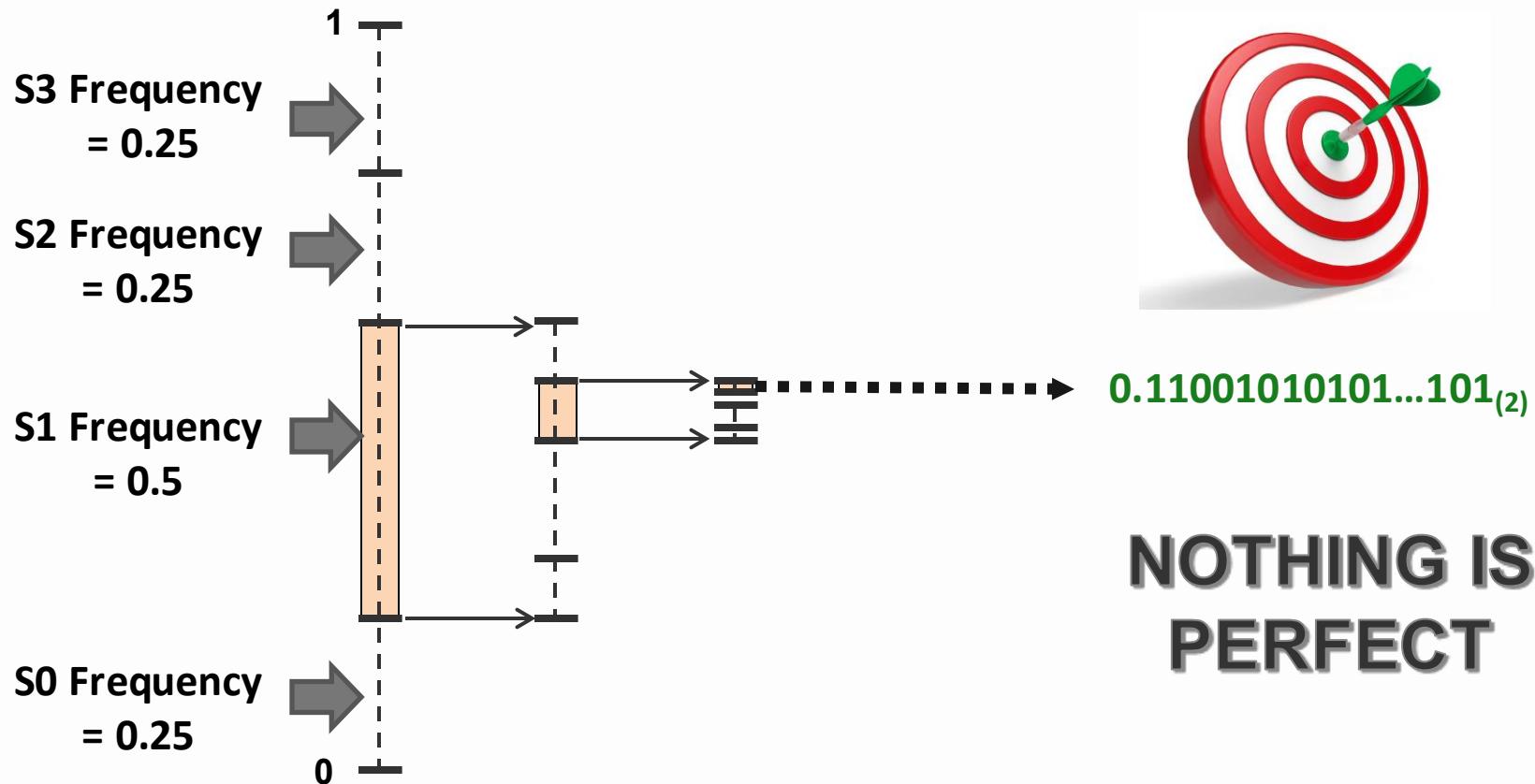


0.1023846489202837462829838393....333292

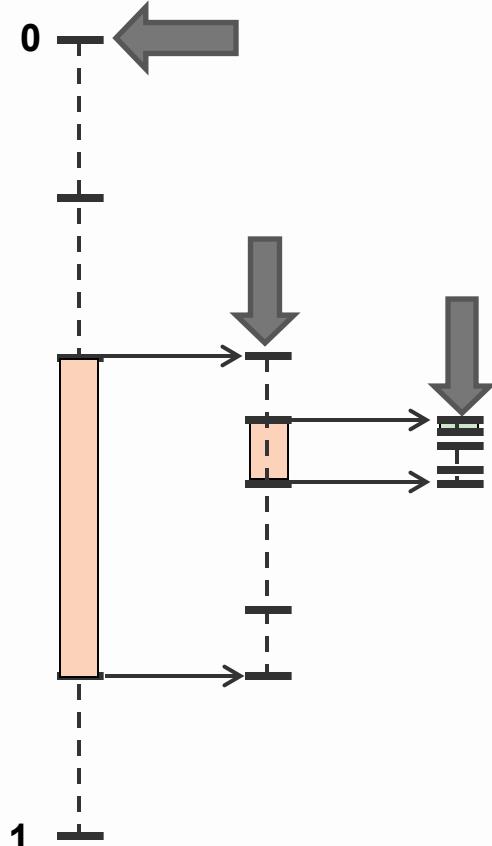
=

0.110101010101010101011110101...111001₍₂₎

Arithmetic Coding – INT2 Example



Arithmetic Coding – Probability Table



Symbol	Frequency
S_1	P_1
S_2	P_2
.	.
.	.
.	.
.	.
.	.
S_N	P_N

INT8
256 Entries
93.6% DRAM Power

Arithmetic Coding – Precision

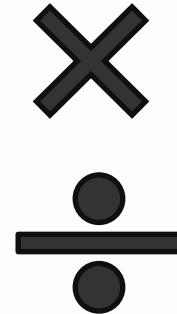


Expandable
Precision

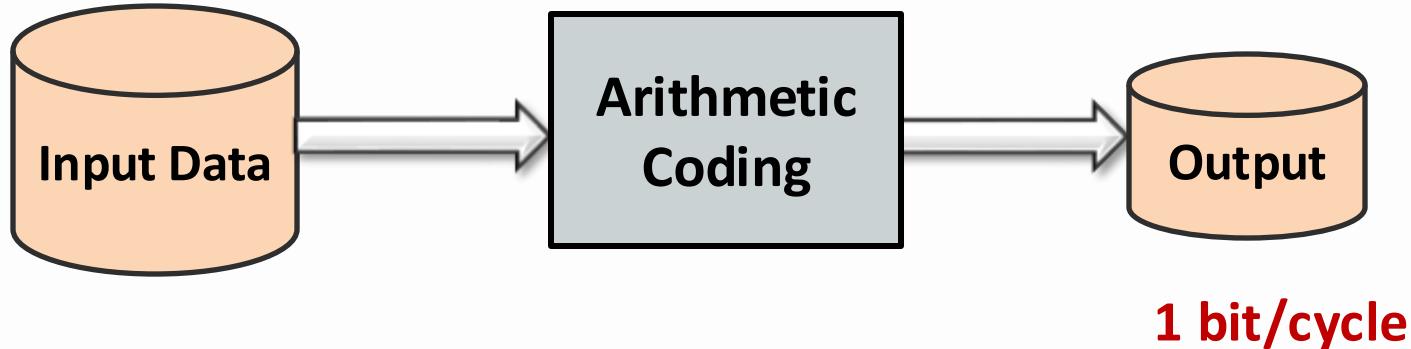


$0.110010.....101_{(2)}$

Costly
Operations



Arithmetic Coding – Bandwidth



Atlanta – Key Idea

VALUE = **BASE** + **OFFSET**

int8 **3 bits** **5 bits**

01010111 = **01000000** + **10111**

256 Entries



↑ **Apply AC**
8 Entries



↑ **Store**

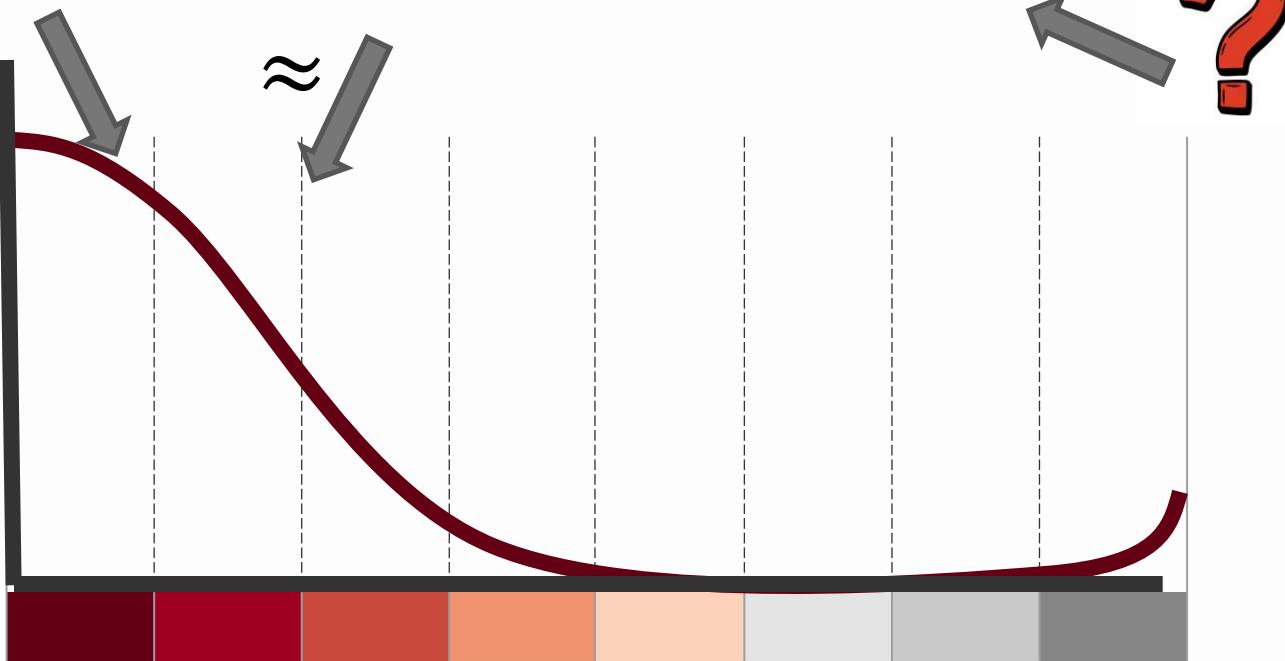
Atlanta – Key Idea

VALUE = **BASE** + **OFFSET**

8 bits

3 bits

5 bits



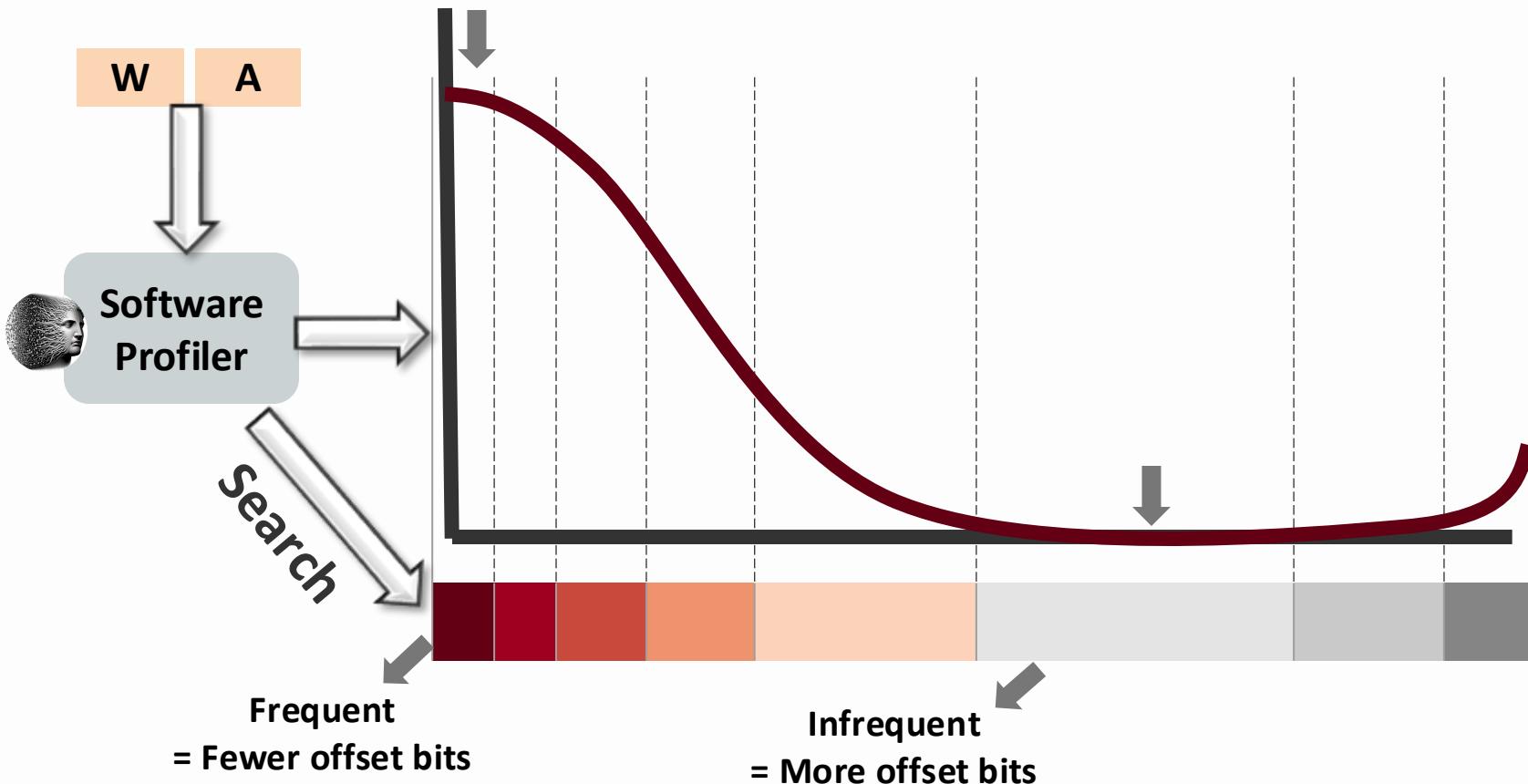


**Software
Profiler**

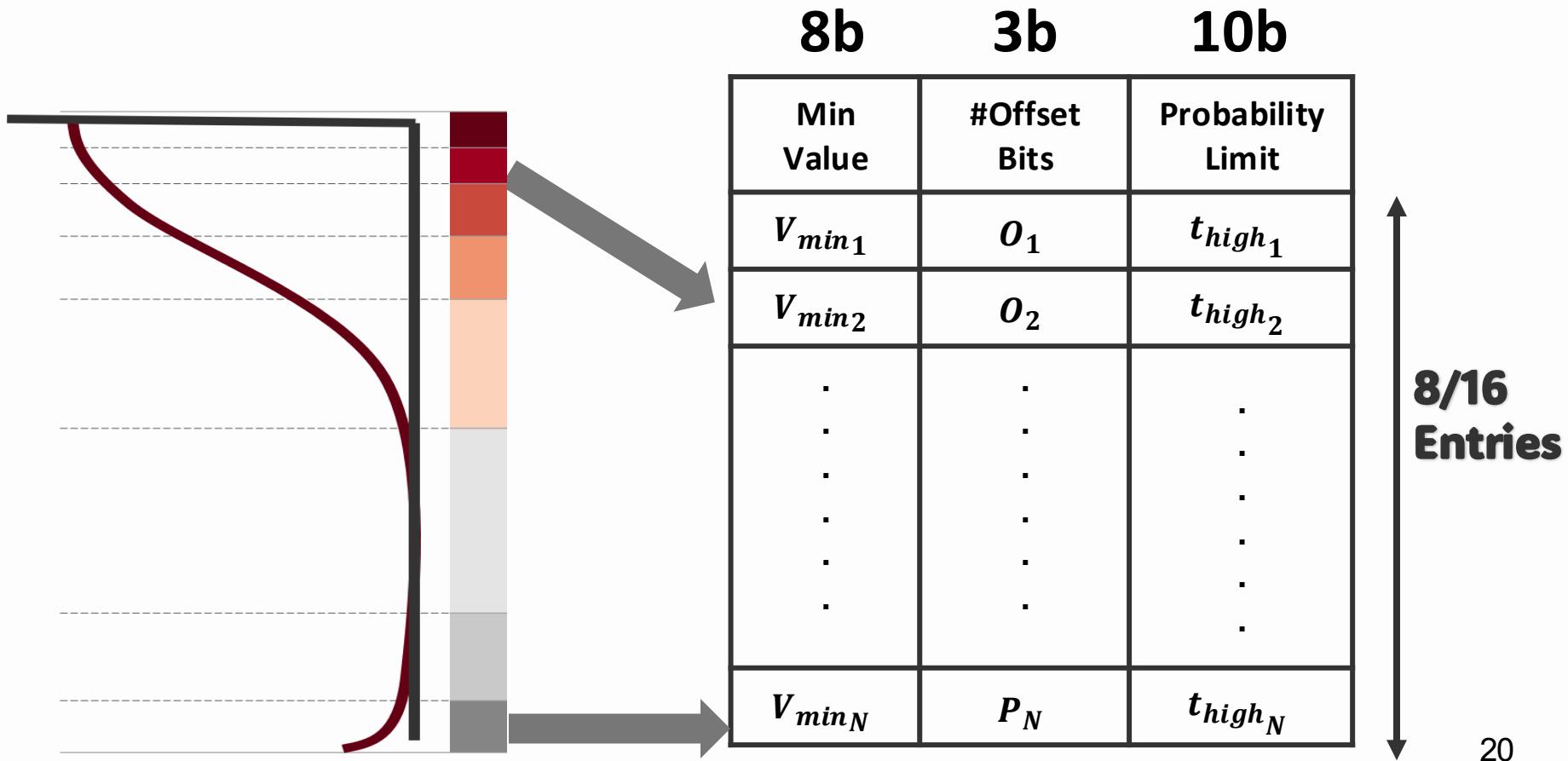


**Hardware Runtime
[Encoder/Decoder]**

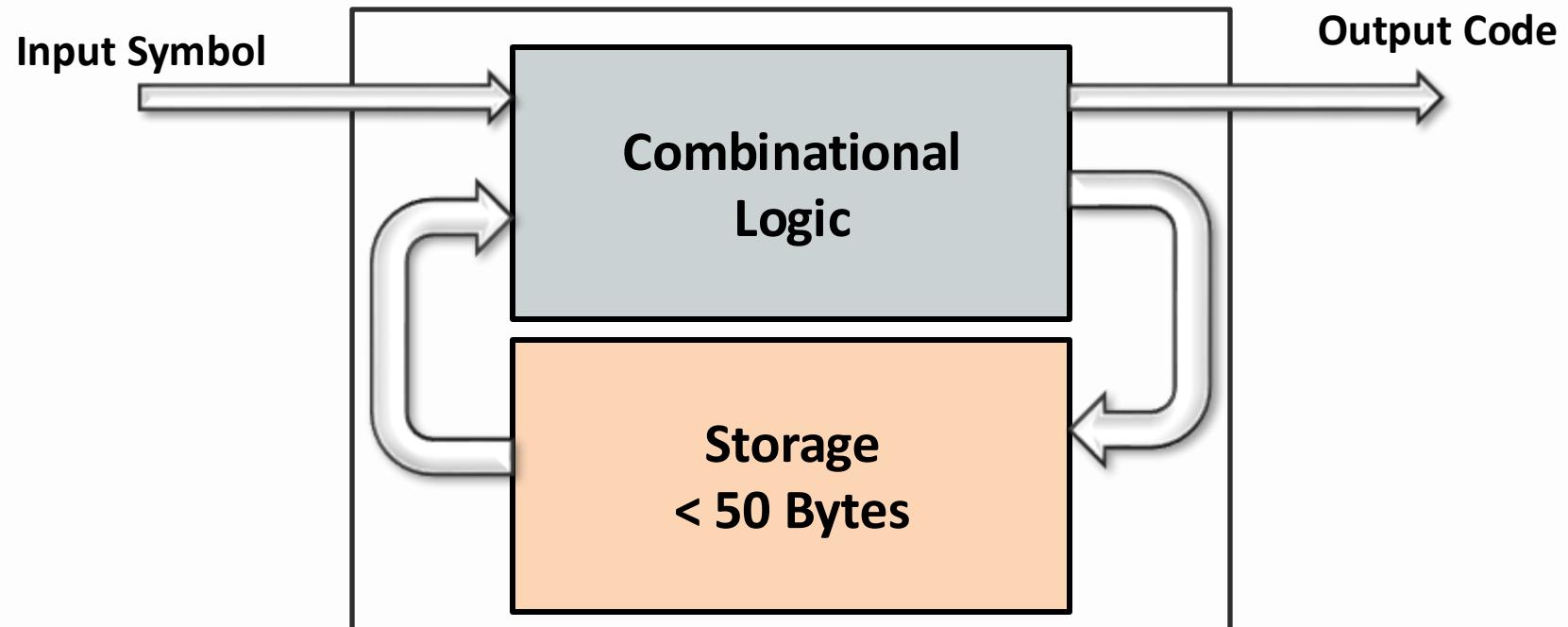
Atalanta – Software Profiler



Atalanta in Software



Atalanta Encoder in Hardware



Atalanta During Inference Evaluation



VS



Bit Plane
Compression
(BPC)

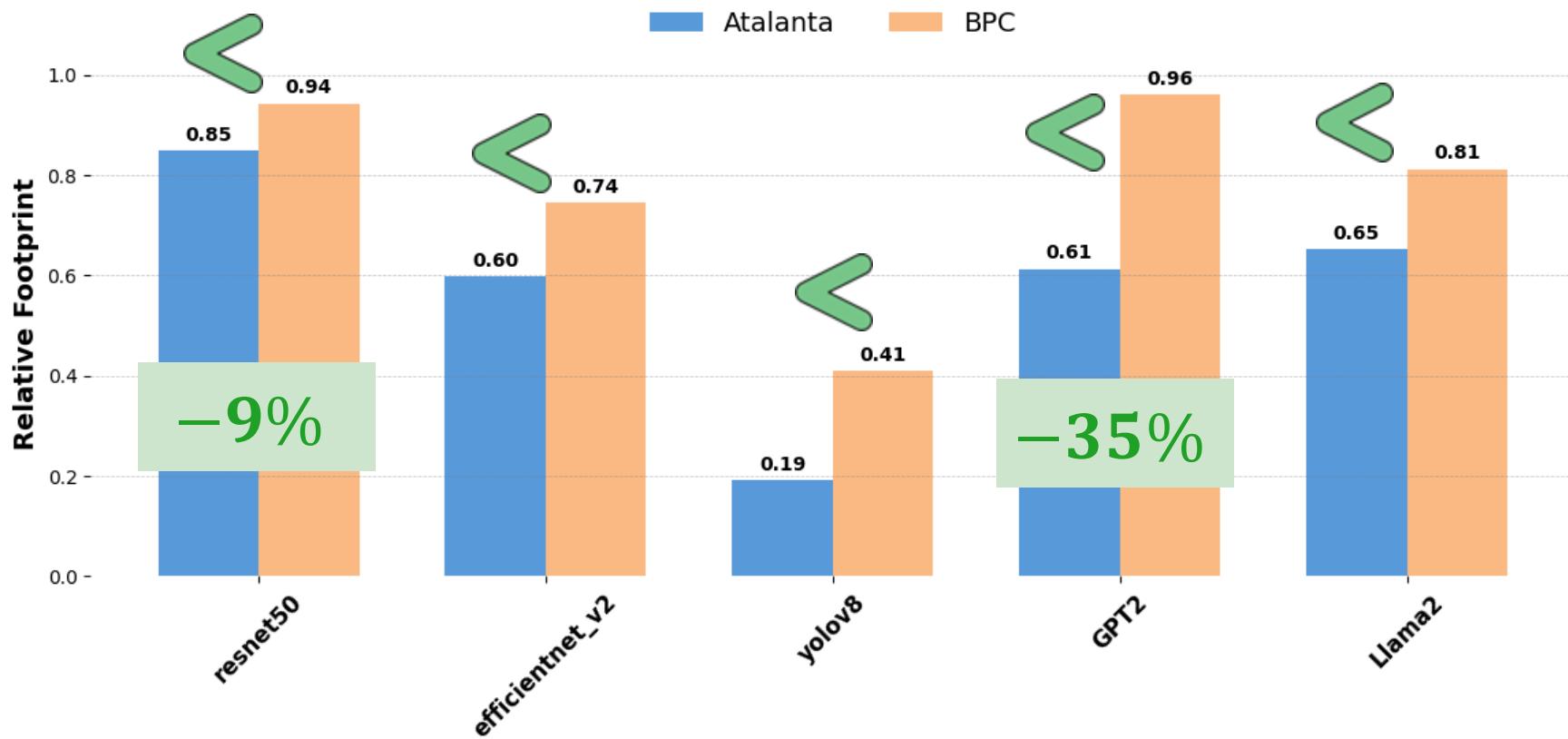


EFFICIENTNET v2

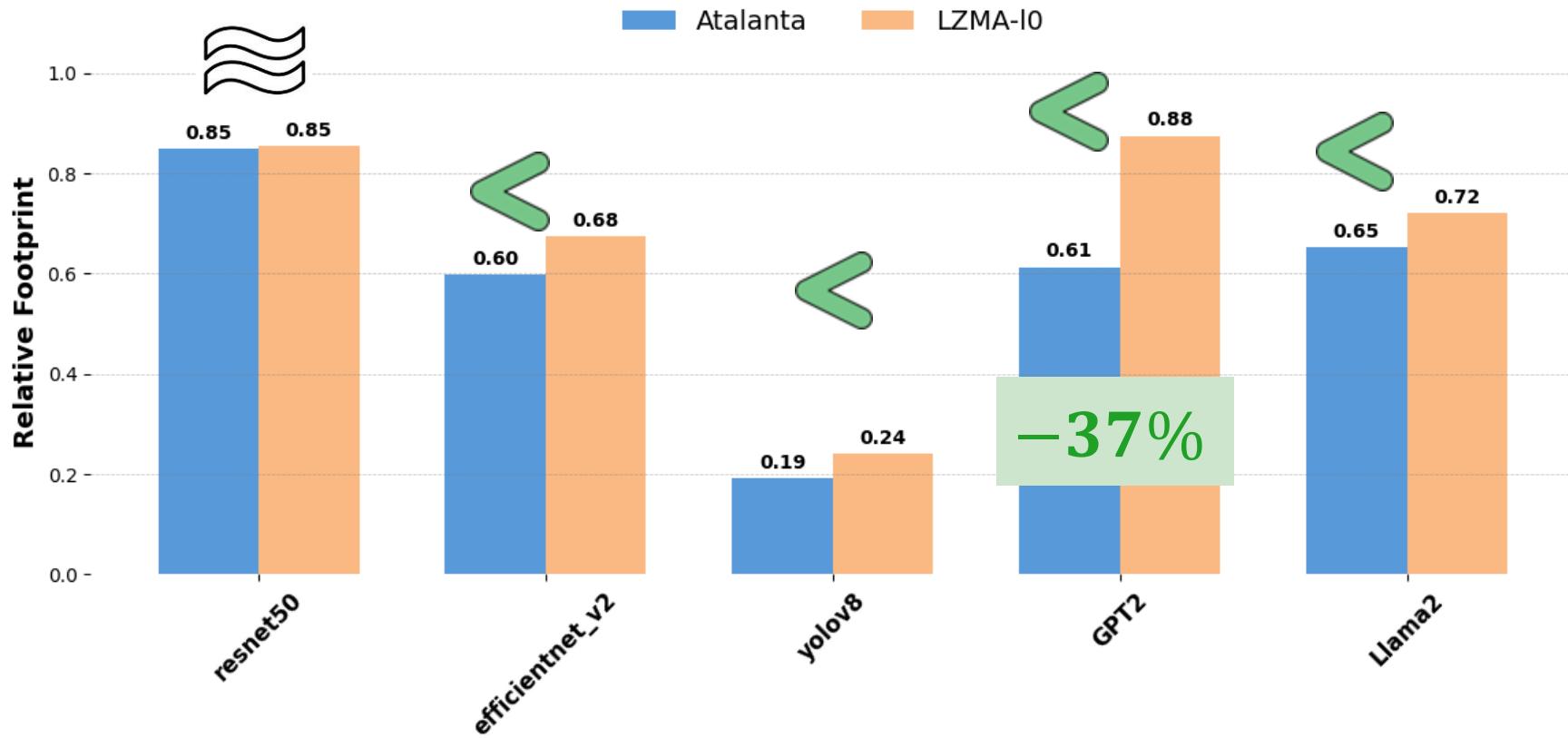


YOLOV8

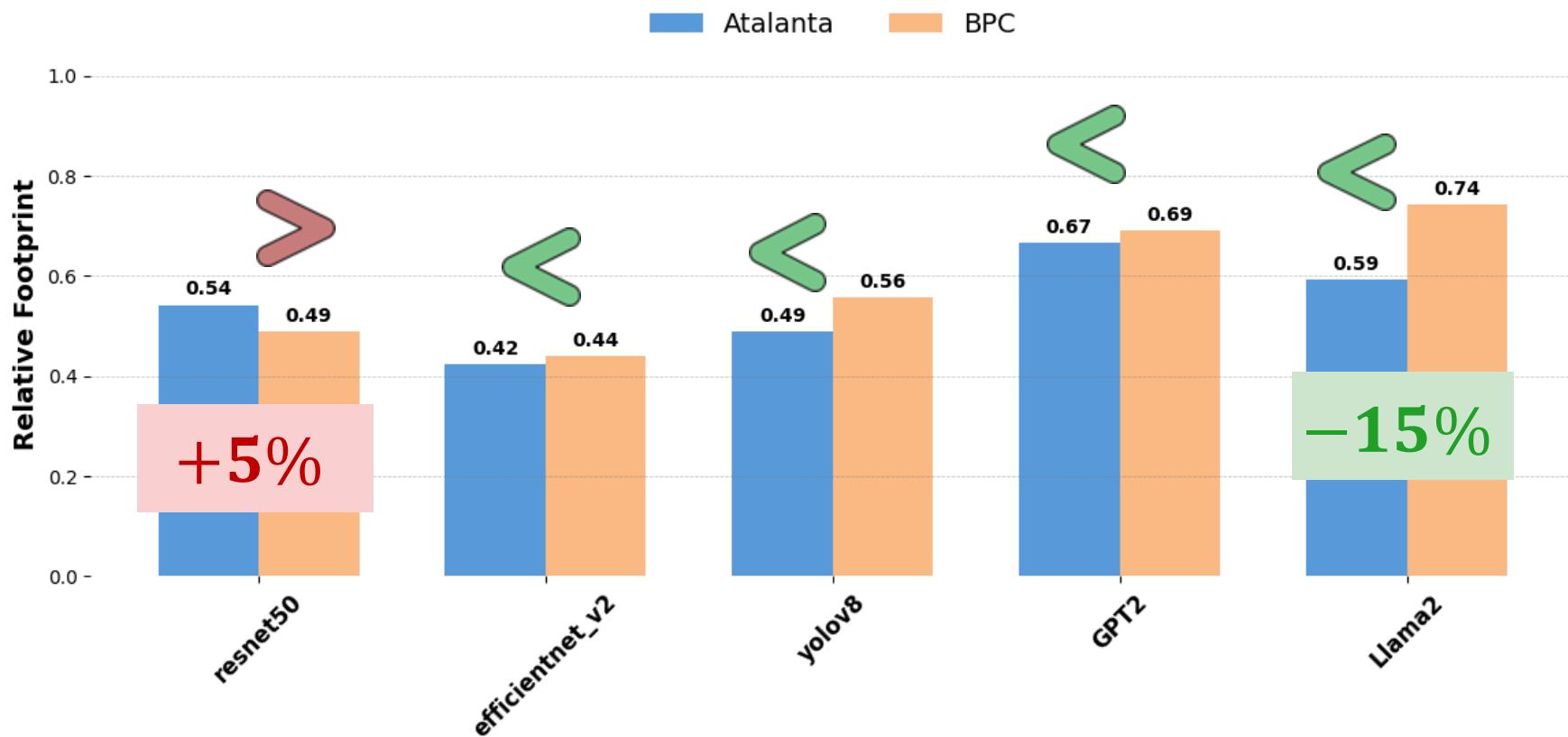
Compression Performance: Weights



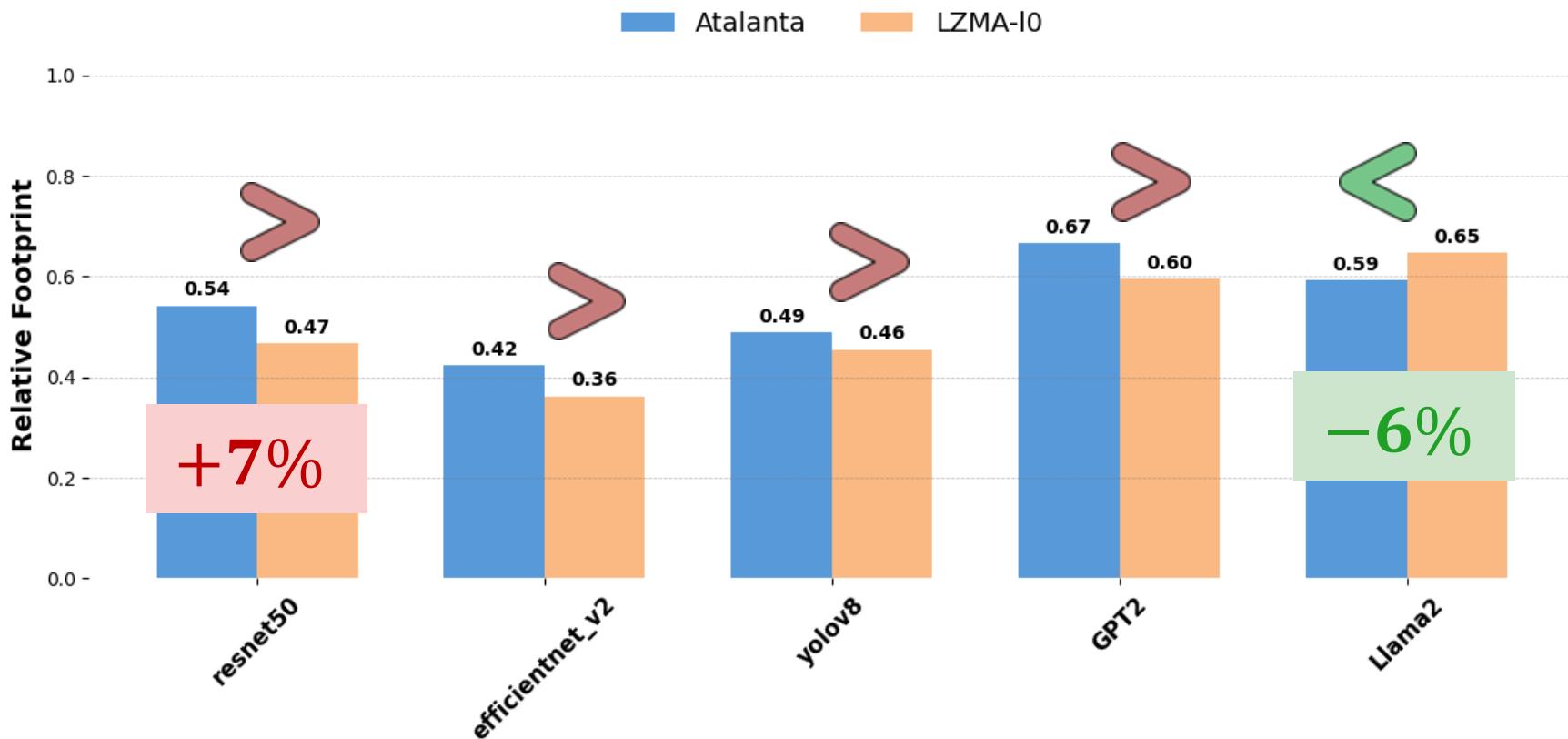
Compression Performance: Weights



Compression Performance: Activations



Compression Performance: Activations



Energy Efficiency and Speedup

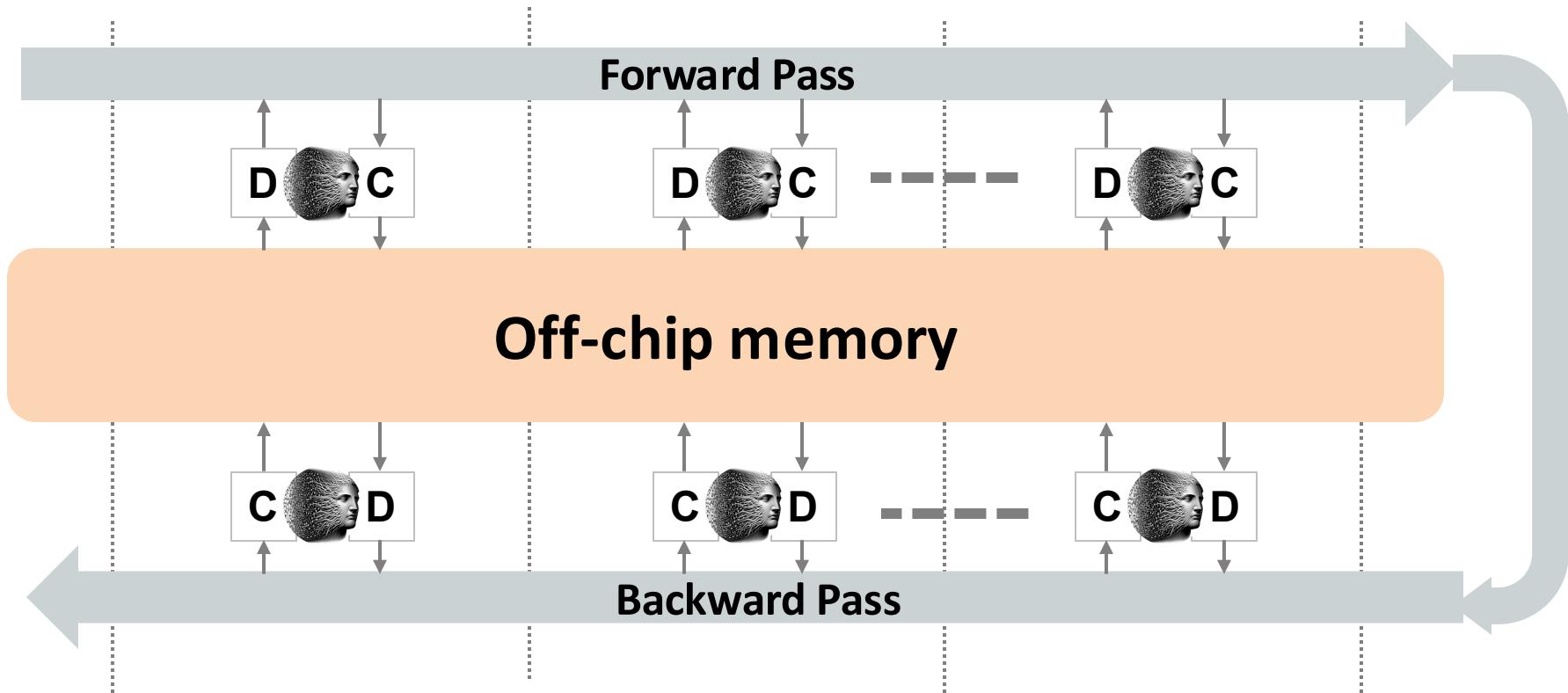
Tensor Cores



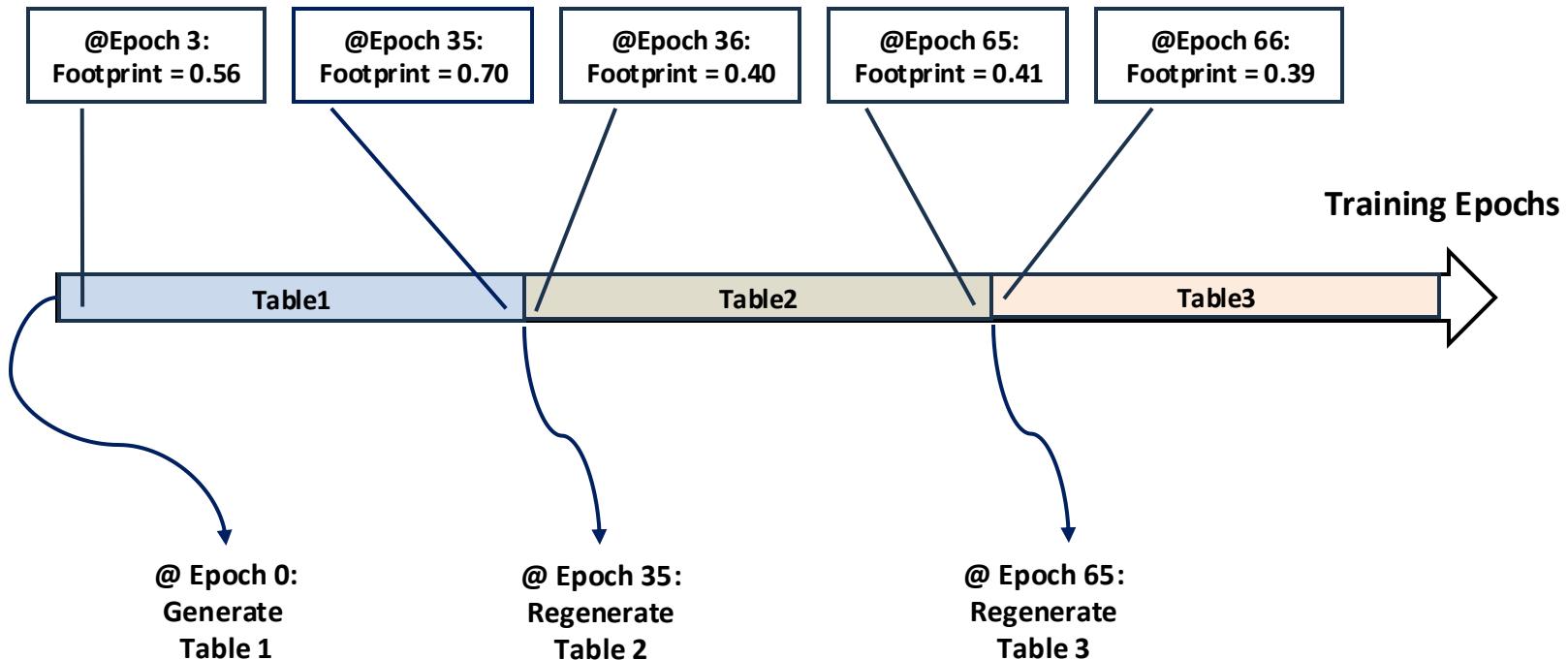
Speedup: 1.44x

Energy Efficiency: 1.37x

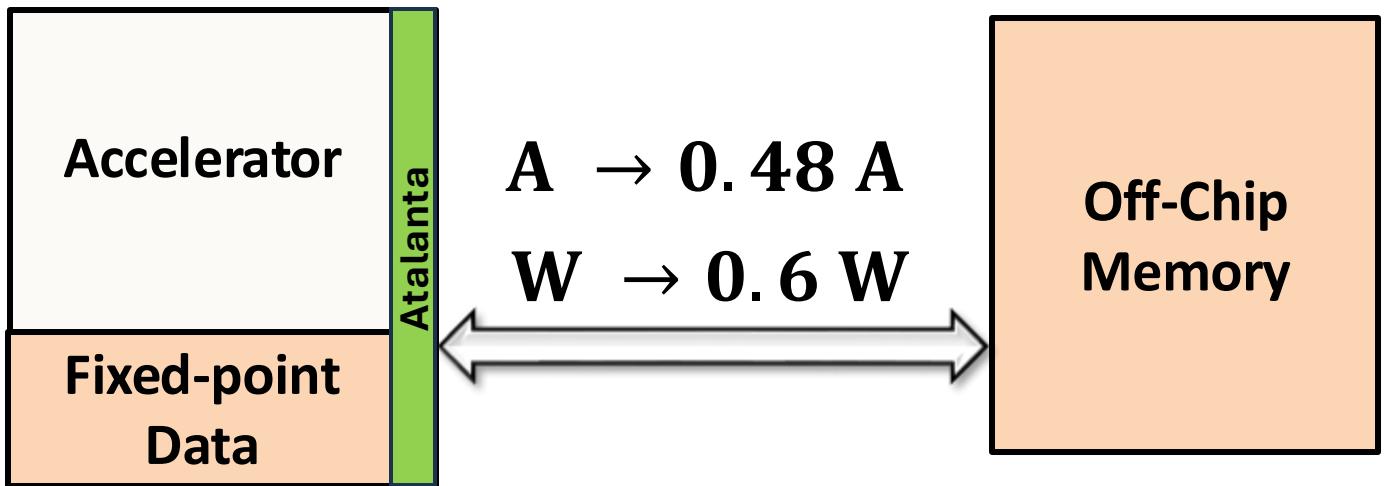
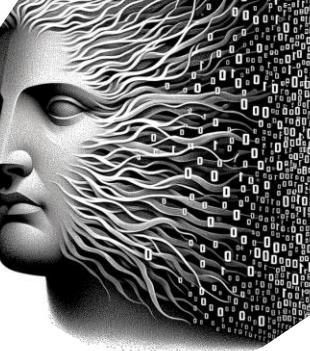
Training for Fixed-Point Inference



Compression Performance During Training



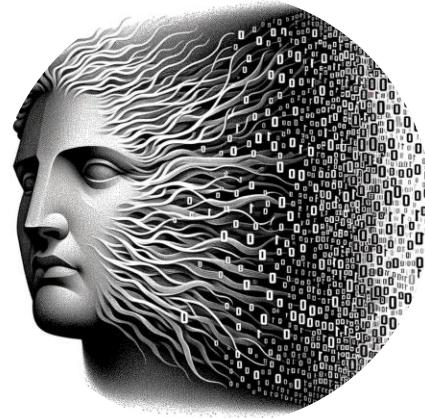
Conclusions



Lossless

Transparent

Low-cost



Thank you!